

A tool for creating and visualizing semantic annotations on relational tables

MAZUMDAR, Suvodeep <<http://orcid.org/0000-0002-0748-7638>> and
ZHANG, Ziqi

Available from Sheffield Hallam University Research Archive (SHURA) at:
<http://shura.shu.ac.uk/16910/>

This document is the author deposited version. You are advised to consult the publisher's version if you wish to cite from it.

Published version

MAZUMDAR, Suvodeep and ZHANG, Ziqi (2016). A tool for creating and visualizing semantic annotations on relational tables. In: GENTILE, Anna Lisa, D'AMATO, Claudia, ZHANG, Ziqi and PAULHEIM, Heiko, (eds.) LD4IE 2016 : Linked data for information extraction : Proceedings of the Fourth International Workshop on Linked Data for Information Extraction co-located with 15th International Semantic Web Conference (ISWC 2016), Kobe, Japan, October 18, 2016. CEUR Workshop Proceedings . CEUR Workshop Proceedings, 2-10.

Copyright and re-use policy

See <http://shura.shu.ac.uk/information.html>

A Tool for Creating and Visualizing Semantic Annotations on Relational Tables

Suvodeep Mazumdar¹ and Ziqi Zhang²

¹Department of Computer Science, University of Sheffield
211 Portobello, Sheffield, UK

²School of Science and Technology, Nottingham Trent University
50 Shakespeare Street, Nottingham, NG1 4FQ

¹s.mazumdar@sheffield.ac.uk, ²ziqi.zhang@ntu.ac.uk

Abstract. Semantically annotating content from relational tables on the Web is a crucial task towards realizing the vision of the Semantic Web. However, there is a lack of open source, user-friendly tools to facilitate this. This paper describes an extension of the TableMiner⁺ system, an open source Semantic Table Interpretation system that automatically annotates Web tables using Linked Data in an effective and efficient approach. It adds a graphical user interface to TableMiner⁺, to facilitate the visualization and correction of automatically generated annotations. This makes TableMiner⁺ an ideal tool for the semi-automatic creation of high-quality semantic annotations on relational tables, which facilitates the publication of Linked Data on the Web.

Keywords: Web table, Named Entity Disambiguation, Semantic Table Interpretation, table annotation, Linked Data

1 Introduction

Recovering semantics from the growing amount of tabular data on the Web is a crucial task in realizing the vision of the Semantic Web. Traditional search engines perform poorly on such data, as they ignore the semantics of tabular structures [2, 3]. Recent years have seen an increase in the research on Semantic Table Interpretation [2, 5, 3, 6], which annotates relational tables using schema and entities defined in a reference knowledge base. The process deals with three types of annotation tasks in tables. Starting with the input of a well-formed relational table, it (1) links entity mentions in content cells to named entities; (2) annotates columns with concepts if they contain entity mentions, or properties of concepts if they contain data literals; and (3) identifies the semantic relations between columns. The annotations created can enable semantic indexing and search of the data, and can be used to create Linked Open Data (LOD).

Semantic Table Interpretation systems are intrinsically difficult to implement, due to, e.g., the complexity of the inter-dependent tasks (e.g., the annotation of a cell depends on that of the containing column and vice versa), and the use of different knowledge bases. TableMiner⁺ [7] is such a method adopting

an incremental, bootstrapping approach that starts by creating preliminary and partial annotations of a table using ‘sample’ data, then using the outcome as ‘seed’ to guide interpretation of remaining contents. This is then followed by a message passing process that iteratively refines results on the entire table to create the final optimal annotations. It has been implemented as open-source software (as part of the STI library¹), however, the system is lacking an intuitive user interface, which has made it difficult to be used by an average person with limited technical knowledge.

This work implements a graphical user interface specifically for TableMiner⁺, to make it an easy-to-use tool for annotating Web tables using Linked Data, and also extend it by enabling users to visualise and correct the generated annotations and Linked Data triples. As a result, data publishers can use TableMiner⁺ for transforming tabular data on the Web into high-quality Linked Data, or creating gold-standard for experiment purposes. The remainder of this paper is structured as follows. Section 2 briefly discusses related work; Section 3 gives an overview of TableMiner⁺; Section 4 introduces the improvement carried out in this work; Section 5 concludes this paper.

2 Related Work

Recent years have seen an increasing number of work on Semantic Table Interpretation. Venetis et al. [4] annotate columns in a table with semantic concepts and identify relations between the subject column (typically containing entities that the table is about) and other columns using a database mined with regular lexico-syntactic patterns such as the Hearst patterns [1]. The database records co-occurrence statistics for each pair of values extracted by such patterns. A maximum likelihood inference model is used to predict the best concepts and relations from candidates using these statistics. Limaye et al. [2] uses a joint inference model, i.e., factor graph to model a table and the interdependencies between its components. Table components are modeled as variables represented as nodes on the graph; then the interdependencies among variables are modeled by factors. The task of inference amounts to searching for an assignment of values to the variables that maximizes the joint probability. Mulwad et al. [3] also uses joint inference with semantic message passing. TableMiner [6] and TableMiner⁺[7] adopt a bootstrapping approach starting by creating preliminary annotations of a table using automatically selected ‘sample’ data in the table, followed by a message passing process that iteratively refines the preliminary annotations to create the final optimal results. These methods differ in terms of the inference models, features and background knowledge bases used. As discussed before, existing tools remains difficult to use due to the lack of a user friendly interface.

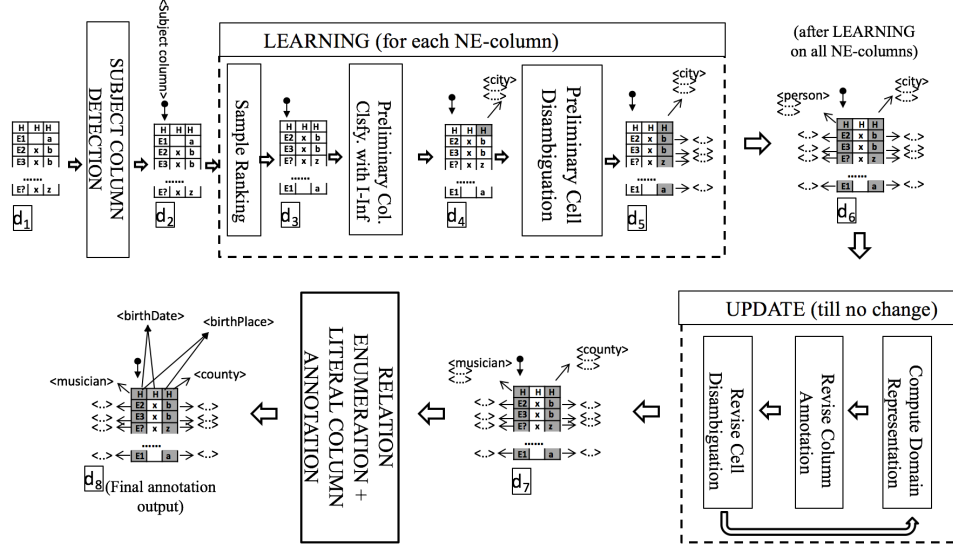
¹ <https://github.com/ziqizhang/sti>

3 Overview of TableMiner⁺

Figure 1 shows a high-level view of the components and workflow of TableMiner⁺. We refer readers to Zhang [7] for details of the methodology. The system can be divided into three major components. Firstly, it detects a ‘subject column’ (**SUBJECT COLUMN DETECTION**), which is the one in the table containing named entities that are subjects of each rows. TableMiner⁺ assumes other columns in a relational table are data describing the subjects. It then identifies other columns that also contain named entities (NE-columns), and performs column classification (assigning a URI from a knowledge base to the column) and cell disambiguation (assigning a URI from a knowledge base to each cell) on these as well as the subject columns. Working with each NE-column at a time, these are further divided into two processes. In the **LEARNING** phase, the system attempts to use a subset (*Sample Ranking*) of rows from the NE-column to infer a concept URI for the column (*Preliminary Col. Classify. with I-Inf*). The idea is that, usually for human-beings, we only need to see some (and rarely do we need to see all) data in a column in order to classify them. However, it is likely that our understanding could be biased because of this ‘partial’ view. And therefore, we call these results ‘preliminary’, which will be optimized later. The LEARNING phase also uses preliminary column annotations as input to guide *Preliminary Cell Disambiguation*. In this part of the process, the assigned concept URI for the column determines the candidate named entities for each row in that column. Next, the preliminary annotations for a column and its content cells are optimized in the **UPDATE** phase. In this phase, the system attempts to ensure annotations on different NE-columns are consistent, e.g., they belong to the same domain (*Compute Domain Representation*). The computation can alter the preliminary annotations in some columns or content cells, which then causes a chain of alterations due to the interdependency of the tasks. A semantic message passing algorithm is implemented to control such update process until convergence. With the column and cell annotations finalized, TableMiner⁺ moves on to infer relations between the subject column and other columns (**RELATION ENUMERATION + LITERAL COLUMN ANNOTATION**). In simple terms, the relation between a subject column and another column is selected based on the relations derived on each row between the pairs of subject entity and data in the other column.

4 Description of the TableMiner⁺ Application Interface

In this section, we describe the TableMiner⁺ user interface and the use of the tool through this interface. We use the implementation distributed as part of the STI library as basis for this work. The STI library provides an implementation of the system introduced in Zhang [7], and a few baseline systems. The library is implemented in Java, and uses DBpedia as the knowledge base. Currently, a Web-based interface consisting of two components are implemented: one that lets users to define, configure and start a table annotation task; and the other that

Fig. 1. Overall component and workflow of TableMiner⁺

lets users to visualise and correct annotation results. In both cases, interaction is achieved via a Web browser².

4.1 Starting a table annotation process

The interface for starting a Semantic Table Annotation task is illustrated in Figure 2³. Users firstly enter the URL of a webpage containing relational tables that are to be annotated. Upon entering the details, the user is shown a preview of the page along with a highlighted list of tables potentially containing relational data. The users then select the tables they wish to annotate. They can also configure the system to alter settings such as feature weights and knowledge base query constraints. The users may provide an email address to subscribe for an automatic alert when the annotation task completes. When the users are satisfied with the configuration and the input, they can click the button to start the task, which will create annotations in JSON format. These will be interpreted and displayed using the visualisation component described below.

4.2 Visualisation and correction of annotations

The JSON files are then passed onto the visualisation component, which consists of two interactive elements: an annotated table and a graph visualisation module.

² However, it is not recommended to deploy TableMiner⁺ as a Web-service as it does not support concurrent access typically found in multi-user environment.

³ Follow <https://github.com/ziqizhang/sti/tree/master/ui> for a demo and on how to use

TableMiner+

Define the extraction task

Q

Enter a URL containing relational tables

https://en.wikipedia.org/wiki/Commedia_all'italiana

Preview

⚙️

Select a reference knowledge base

(Disabled. Only DBpedia is supported.)

DBpedia (SPARQL endpoint, V3.9 ontology)

⚙️

Select a table parser

(To identify relational tables from input files.)

Generic (extracts all n x m tables with horizontal headers)

✉️

Enter an Email address for notification

⚙️

Advanced settings

Advanced users only. You may change settings such as feature weights, knowledge base query constraints, NLP resources used by the system, etc.

Reset

Start

Romanzo popolare (Popular Romance) and *Amici miei*.

Other notable film makers of the genre were Pasquale Festa Campanile, Ettore Scola, Luigi Comencini, Steno, Antonio Pietrangeli, Nanni Loy and Lina Wertmüller.

And the script writers Age & Scarpelli, Leo Benvenuti, Piero De Bernardi, Rodolfo Sonego, Suso Cecchi d'Amico, Sergio Amidei

Notable films

☒

Notable films (check this box to select this table for annotation)

Title	Year	Director	Notable actors
<i>Big Deal on Madonna Street</i>	1958	Mario Monicelli	Marcello Mastroianni, Vittorio Gassman, Totò
<i>The Great War</i>	1959	Mario Monicelli	Vittorio Gassman, Alberto Sordi, Silvana Mangano
<i>Il vedovo</i>	1959	Dino Risi	Alberto Sordi, Franca Valeri
<i>Love and Larceny</i>	1959	Dino Risi	Vittorio Gassman
<i>Everybody Go Home</i>	1960	Luigi Comencini	Alberto Sordi
<i>Adua e le compagne</i>	1960	Antonio Pietrangeli	Simone Signoret, Marcello Mastroianni, Sandra Milo
<i>A Difficult Life</i>	1961	Dino Risi	Alberto Sordi
<i>Audace colpo dei soliti ignoti</i>	1961	Nanni Loy	Vittorio Gassman, Claudia Cardinale, Nino Manfredi
<i>The Fascist</i>	1961	Luciano Salce	Ugo Tognazzi, Stefania Sandrelli
<i>Divorce, Italian Style</i>	1962	Pietro Germi	Marcello Mastroianni, Stefania Sandrelli
<i>Boccaccio '70</i>	1962	Mario Monicelli, Federico Fellini, Luchino Visconti, Vittorio De Sica	

Fig. 2. User input interface, where the user enters source URL containing the respective relational tables. The preview button provides an idea of the table that will be annotated, and proceeding with start initiates the extraction process. An email address can be provided if the user prefers not to wait and be notified once the extraction process has completed.

The annotated table is the first point of interaction with the user, and presents the original table, annotated with the entities, concepts and relations identified by TableMiner⁺. The first step for the UI is to investigate the header cells of the table - TableMiner⁺ creates a set of candidate concepts that best describe the header and the data in the column. Each associated concept has a score indicating the system's confidence. This set of candidate concepts is presented as a dropdown with the scores (Figure 3 section B). Users can select any of the concepts to indicate a more appropriate annotation by clicking on the respective concept. Concepts are further encoded on the basis of scores (the highest scoring concepts are indicated in green, while the lowest in red), which provides an indication of the confidence in content cell annotations can be visualised in the same way.

As can also be seen from the figure, some entities have already been recognized, while some haven't. In such cases the user can provide a URI that is appropriate for any missing annotation, this can be done by clicking the relevant cell, which will provide a prompt for a text input (Figure 4). Further SPARQL queries can also be triggered to the respective endpoints (based on the user customisations) that can identify any missing annotations.

While tables can provide a clean annotated replication of the original source document, with the added ability for users to provide their annotations and

C	Title	Year	Director	Notable actors
	Italian Films 1.66		Film Director 3.7	Actor 2.88
Q	Big Deal on Madonna Street http://dbpedia.org/resource/Big_Deal_on_Madonna_Street , 2.06	1958	Film Maker 2.91 Yago Legal Actor 2.61	Marcello Mastroianni http://dbpedia.org/resource/Marcello_Mastroianni , 2.39
Q	The Great War http://dbpedia.org/resource/The_Great_War , 1.60	1959	Organism 2.38 Producer 2.38	Vittorio Gassman http://dbpedia.org/resource/Vittorio_Gassman , 2.19
Q	Il vedovo http://dbpedia.org/resource/Il_vedovo , 2.56	1959	Physical Entity 2.38 Natural Person 2.38	Alberto Sordi http://dbpedia.org/resource/Alberto_Sordi , 2.25
Q	Love and Larceny	1959	Whole 2.38	Vittorio Gassman http://dbpedia.org/resource/Vittorio_Gassman , 2.19
Q	Everybody Go Home http://dbpedia.org/resource/Everybody_Go_Home , 2.27	1960	Creator 2.38 Italian Atheists 2.06	Alberto Sordi http://dbpedia.org/resource/Alberto_Sordi , 2.25
Q	Adua e le compagne	1960	Italian Screenwriters 1.58 Film Directors Who Committed Suicide 1.09	Simone Signoret http://dbpedia.org/resource/Simone_Signoret , 2.05
Q	A Difficult Life http://dbpedia.org/resource/A_Difficult_Life , 2.32	1961	Screenwriter 1.08 People From Viareggio 0.66	Alberto Sordi http://dbpedia.org/resource/Alberto_Sordi , 2.25
Q	Audace colpo dei soliti ignoti http://dbpedia.org/resource/Audace_colpo_dei_soliti_ignoti , 1.91	1961	People From Sal C 0.25 People From Sal C 0.25	Vittorio Gassman http://dbpedia.org/resource/Vittorio_Gassman , 2.19
Q	The Fascist http://dbpedia.org/resource/The_Fascist , 2.53	1961	Luciano Salce http://dbpedia.org/resource/Luciano_Salce , 2.16	Ugo Tognazzi http://dbpedia.org/resource/Ugo_Tognazzi , 2.49
Q	Divorce, Italian Style http://dbpedia.org/resource/Divorce_Italian_Style , 1.86	1962	Pietro Germi http://dbpedia.org/resource/Pietro_Germi , 2.19	Marcello Mastroianni http://dbpedia.org/resource/Marcello_Mastroianni , 2.39
Q	Boccaccio '70	1962	Mario Monicelli http://dbpedia.org/resource/Mario_Monicelli , 2.60	
Q	The Easy Life http://dbpedia.org/resource/The_Easy_Life , 1.97	1962	Dino Risi http://dbpedia.org/resource/Dino_Risi , 2.73	Vittorio Gassman http://dbpedia.org/resource/Vittorio_Gassman , 2.19
Q	The Last Judgement	1962	Vittorio De Sica http://dbpedia.org/resource/Vittorio_De_Sica , 2.43	Vittorio Gassman http://dbpedia.org/resource/Vittorio_Gassman , 2.19
Q	Malioso	1962	Alberto Lattuada http://dbpedia.org/resource/Alberto_Lattuada , 2.46	Alberto Sordi http://dbpedia.org/resource/Alberto_Sordi , 2.25
Q	March on Rome	1962	Dino Risi http://dbpedia.org/resource/Dino_Risi , 2.73	Vittorio Gassman http://dbpedia.org/resource/Vittorio_Gassman , 2.19
Q	The Conjugal Bed	1963	Marco Ferreri http://dbpedia.org/resource/Marco_Ferri , 2.56	Ugo Tognazzi http://dbpedia.org/resource/Ugo_Tognazzi , 2.49
Q	I mostri http://dbpedia.org/resource/I_mostri , 2.54	1963	Dino Risi http://dbpedia.org/resource/Dino_Risi , 2.73	Vittorio Gassman http://dbpedia.org/resource/Vittorio_Gassman , 2.19
Q	Alta infedeltà	1965	Mario Monicelli http://dbpedia.org/resource/Mario_Monicelli , 2.60	Ugo Tognazzi http://dbpedia.org/resource/Ugo_Tognazzi , 2.49
Q	Il diavolo http://dbpedia.org/resource/Il_diavolo , 2.54	1963	Gian Luigi Polidoro http://dbpedia.org/resource/Gian_Luigi_Polidoro , 2.12	Alberto Sordi http://dbpedia.org/resource/Alberto_Sordi , 2.25
Q	Il Boom http://dbpedia.org/resource/Il_Boom , 2.67	1963	Vittorio De Sica http://dbpedia.org/resource/Vittorio_De_Sica , 2.43	Alberto Sordi http://dbpedia.org/resource/Alberto_Sordi , 2.25
Q	Yesterday, Today and Tomorrow	1963	Vittorio De Sica http://dbpedia.org/resource/Vittorio_De_Sica , 2.43	Marcello Mastroianni http://dbpedia.org/resource/Marcello_Mastroianni , 2.39

Fig. 3. TableMiner UI Annotated Table for the page available at https://en.wikipedia.org/wiki/Commedia_all%27italiana

Q	Divorce, Italian Style http://dbpedia.org/resource/Divorce_Italian_Style , 1.86	Q	Divorce, Italian Style http://dbpedia.org/resource/Divorce_Italian_Style , 1.86
Q	Boccaccio '70	Q	Boccaccio '70 https://en.wikipedia.org/wiki/Boccaccio_%2770
Q	The Easy Life http://dbpedia.org/resource/The_Easy_Life , 1.97	Q	The Easy Life http://dbpedia.org/resource/The_Easy_Life , 1.97

Fig. 4. Editing cell contents to add annotations - for missing/wrong annotations (as seen in the figure left), clicking on the relevant cell will enable editing (right). All changes made on the tables are stored until finally pushed to the backend database

correct any mistakes they can observe, a further need may arise for greater customisation and control of annotations. This provides users with means to visualise (and annotate) possible relations among table columns, in addition to visualising possible candidate annotations. The next aspect of the UI is the graph visualisation, which is invoked from the ‘inspect’ button on the first cell of each table row (Figure 3, Section C). As an example, the header and it’s relevant candidate concepts have been plotted as a graph in Figure 5.

Header cells are shown as nodes labelled with the header columns (0-3), while the candidate classes are shown as nodes, linked with header elements. The most relevant class is shown with a strong link, while the others are presented as dashed lines. Clicking on an individual node makes all other nodes more transparent, and hence keeps the current node and link in focus. Right-clicking the dashed ones annotate the relevant header cell with the respective concept, which will then confirm the change with a strong link (here, a straight thick line). Header cells are also linked with each other with dashed lines, which is

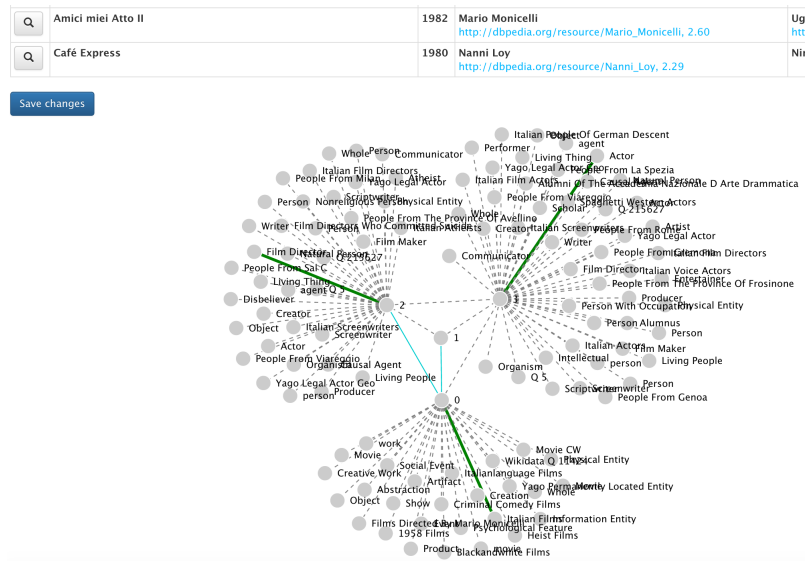


Fig. 5. Visualising table rows as a node-link graph

interpreted as only an indicative relation. However, if TableMiner⁺ creates any relations between the columns, it is reflected as straight lines as can be seen in Figure 5.

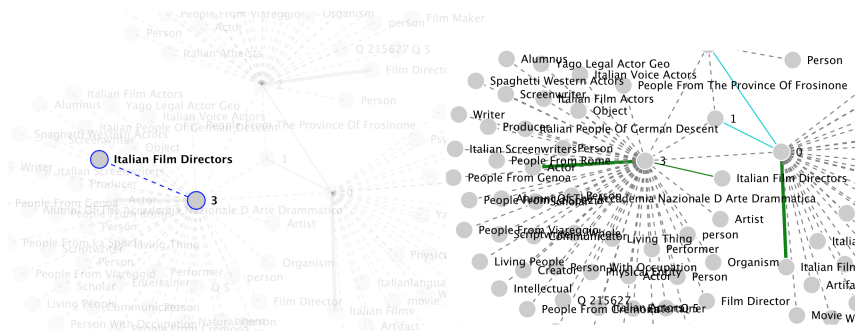


Fig. 6. Clicking on nodes focusses the view on the clicked node and immediate connections to other nodes. Right clicking a node creates a stronger connection (right, a new link is established with 'Italian Film Directors') which is interpreted as another annotation for the column in the original table

Described Scenario In the example shown so far, the source URL (https://en.wikipedia.org/wiki/Commedia_all'italiana) describes movies re-

leased that belong to an Italian film genre. The extraction process in TableMiner⁺ annotated several cells of the table selected by the user in Figure 2 (Notable films), however several films could not be identified. This is made evident when the user visualises the annotated table (Figure 3). Missing cells can then be manually annotated by adding URLs if the user can provide any unidentified ones (Figure 4). For example, the user observes the cell ‘Boccaccio 70’ could not be identified and hence chose to manually add the resource URL. Further inspecting the different concepts, users can click on the ‘inspect’ button to visualise the concepts on a node-link graph (Figure 5). While interacting with the graph, the user notes that since the original webpage discussed Italian movie genre, the ‘Italian Film Directors’ concept would be appropriate to describe the first column in the table. Hence the user can right-click on the concept to add a new annotation for the column. Each row in the table can be visualised as a graph, and hence the user can introduce row-specific annotations as well. Finally, when all annotations are completed, the user can click on ‘Save changes’ to submit all annotations to TableMiner⁺.

5 Conclusion

This paper introduced a graphical user interface for TableMiner⁺ to facilitate the semi-automatic creation of high quality Linked Data and annotations on Web tables. Future work will extend the system to support, e.g., different knowledge bases, other algorithms, fine-grained task definition that enable batch processing and zoning on tables (e.g., specific columns). Furthermore, we will also explore different visualisations and mechanisms for users to introduce new annotations, visualising relevant sections of ontologies while exploring table annotations. We also have a series of user evaluations planned to understand how users can make use of the user interface.

Acknowledgement This work is funded by the EU FP7 WeSenseIt (grant agreement 308429)⁴ and EU Horizon 2020 Seta (grant agreement 688082)⁵ projects. We also thank the ADEQUATE⁶ project team under the lead of Dr Tomas Knap for contributing valuable design ideas.

References

1. Marti A. Hearst. Automatic acquisition of hyponyms from large text corpora. In *Proceedings of the 14th Conference on Computational Linguistics - Volume 2*, COLING ’92, pages 539–545, Stroudsburg, PA, USA, 1992. Association for Computational Linguistics.

⁴ <http://wesenseit.eu/>

⁵ <http://setamobility.eu/>

⁶ <http://www.adequate.at>

2. Girija Limaye, Sunita Sarawagi, and Soumen Chakrabarti. Annotating and searching web tables using entities, types and relationships. *Proceedings of the VLDB Endowment*, 3(1-2):1338–1347, 2010.
3. Varish Mulwad, Tim Finin, and Anupam Joshi. Semantic message passing for generating linked data from tables. In *International Semantic Web Conference (1)*, Lecture Notes in Computer Science, pages 363–378. Springer, 2013.
4. Petros Venetis, Alon Halevy, Jayant Madhavan, Marius Pasca, Warren Shen, Fei Wu, Gengxin Miao, and Chung Wu. Recovering semantics of tables on the web. *Proceedings of VLDB Endowment*, 4(9):528–538, June 2011.
5. Jingjing Wang, Haixun Wang, Zhongyuan Wang, and Kenny Q. Zhu. Understanding tables on the web. In *Proceedings of the 31st international conference on Conceptual Modeling*, ER’12, pages 141–155, Berlin, Heidelberg, 2012. Springer-Verlag.
6. Ziqi Zhang. Towards effective and efficient semantic table interpretation. In *Proceedings of the 13th International Semantic Web Conference*, pages 487–502, 2014.
7. Ziqi Zhang. Effective and efficient semantic table interpretation using tableminer+. *Semantic Web Journal*, Accepted. Tracking: 1339-2551, 2016.